**Research Article**

# Fine-Tuning or Prompting? The Real Economics of ChatGPT in Enterprises

[1]Chauhan Minal Vinodbhai, [2]Bhura Parul Abhesinh, [3]Chowdhary Nehal Sureshbhai

[1]*Assistant Professor*
*minal19ch@gmail.com*
[2]*Assistant Professor*
*bhura2parul@gmail.com*
[3]*Assistant Professor*
*nehal.chawdhary12@gmail.com*

| ARTICLE INFO | ABSTRACT |
|---|---|
| | The rapid adoption of large language models (LLMs) such as ChatGPT has generated significant debate within enterprises over whether to rely on prompt engineering or invest in fine-tuning for domain-specific applications. While prompting offers flexibility and low entry costs, fine-tuning provides potential gains in task specialization, consistency, and efficiency. This paper investigates the comparative economics of fine-tuning versus prompting within organizational settings, with a focus on deployment costs, performance trade-offs, and strategic scalability. Using a mixed-methods framework that combines cost modeling, benchmark analysis, and enterprise case studies, we evaluate scenarios across customer support, compliance documentation, and knowledge management. Findings reveal that prompting remains more cost-effective for exploratory and low-volume use cases, whereas fine-tuning becomes economically viable at scale when consistency and compliance are prioritized. The study highlights trade-offs between short-term experimentation and long-term operational integration, offering a structured decision-making framework for enterprise leaders. By situating the analysis within broader debates on responsible AI adoption, this paper contributes to both scholarly and managerial understanding of how enterprises can maximize returns on LLM investments.<br><br>**Keywords:** ChatGPT, fine-tuning, prompt engineering, enterprise AI, cost analysis |

## Introduction

Large language models (LLMs) such as OpenAI's ChatGPT are increasingly being integrated into enterprise operations to streamline workflows, reduce costs, and enhance decision-making (Dwivedi et al., 2023). From automating customer service to drafting compliance reports, enterprises are exploring how to deploy LLMs most effectively. Yet a key strategic question persists: should organizations rely primarily on prompt engineering, where task instructions are optimized in natural language, or invest in fine-tuning, where the base model is retrained on domain-specific data? This decision is not merely technical; it has profound economic and organizational implications, influencing scalability, data governance, compliance risk, and return on investment (ROI).

Prompt engineering has gained traction because it enables non-technical staff to elicit high-quality outputs without specialized training data (White et al., 2023). It supports rapid experimentation and flexible adaptation across diverse tasks. However, prompting may suffer from inconsistency, requiring continuous human oversight, and may not achieve the compliance standards demanded in regulated industries (Bommasani et al., 2022).

In contrast, fine-tuning involves adapting the LLM weights using proprietary or curated datasets. While fine-tuning offers greater task fidelity and consistency, it comes with higher upfront costs, demands on data infrastructure, and ongoing retraining as enterprise needs evolve (Hu et al., 2023). For many organizations, the trade-off centers on whether the gains in efficiency and compliance justify the higher capital expenditure.

Despite the surge of interest in LLM deployment, existing scholarship has primarily focused on technical comparisons of prompting and fine-tuning, with limited attention to their economic implications for enterprises. Current studies

**Research Article**

often present performance benchmarks (Wang et al., 2023) but seldom examine cost structures, scalability thresholds, and long-term value capture. Moreover, there is a lack of comprehensive frameworks guiding decision-makers in aligning LLM strategies with organizational goals, industry regulation, and resource constraints.

This paper addresses these gaps by systematically analyzing the economics of fine-tuning versus prompting in enterprise contexts. Specifically, we seek to answer the following research questions:

**RQ1:** Under what circumstances is fine-tuning more cost-effective than prompting for enterprises?

**RQ2:** How do task characteristics (e.g., compliance sensitivity, volume, and complexity) shape the economic trade-offs between prompting and fine-tuning?

**RQ3:** What decision-making framework can enterprises adopt to balance short-term flexibility with long-term efficiency?

This study contributes to academic and managerial debates in three key ways. First, it develops a comparative cost–benefit framework for evaluating prompting and fine-tuning strategies. Second, it provides empirical insights by modeling enterprise-scale use cases across customer service, compliance, and knowledge management. Third, it situates the findings within broader concerns of responsible AI adoption, highlighting issues of governance, ethics, and workforce implications.

The paper is organized as follows: Section 2 reviews the literature on prompting, fine-tuning, and enterprise AI economics. Section 3 outlines the research methodology, including cost modeling and case study design. Section 4 presents results from empirical analyses and simulations. Section 5 discusses implications, trade-offs, and limitations. Section 6 concludes with a decision-making framework and directions for future research.

## Literature Review

### 2.1 Large Language Models in Enterprise Adoption

The emergence of large language models (LLMs) has transformed the enterprise AI landscape by enabling generalized capabilities such as reasoning, summarization, and knowledge synthesis without task-specific training (Bommasani et al., 2022). ChatGPT, among the most widely deployed, has been adopted in domains such as customer support, financial services, healthcare, and legal compliance (Dwivedi et al., 2023). Enterprises are attracted to LLMs not only for productivity gains but also for their potential to reduce costs of routine knowledge-intensive tasks (Hu et al., 2023). However, their adoption presents a strategic dilemma: whether to rely on prompting techniques that adapt the model at inference or invest in fine-tuning to embed enterprise-specific knowledge.

### 2.2 Prompt Engineering: Flexibility and Low-Cost Adaptation

Prompt engineering is increasingly recognized as a practical, low-barrier method for harnessing LLMs in enterprise contexts. Research suggests that well-crafted prompts can significantly improve LLM accuracy across reasoning and classification tasks without modifying model parameters (White et al., 2023). This approach allows rapid prototyping and experimentation, particularly valuable for organizations lacking large proprietary datasets.

However, reliance on prompting poses challenges. First, prompts often produce inconsistent results, with outputs sensitive to phrasing variations (Liu et al., 2023). Second, prompt-based methods may struggle with compliance-sensitive contexts where consistent decision-making is legally mandated, such as financial auditing or legal drafting (Gilardi et al., 2023). Finally, prompts can be difficult to scale operationally since maintaining prompt libraries across diverse tasks requires ongoing expert oversight.

Despite these limitations, studies show that prompting remains economically favorable in exploratory settings, where cost minimization and agility outweigh concerns of reliability (Zhang et al., 2023). This explains its popularity among enterprises experimenting with AI but not yet committed to full-scale integration.

**Research Article**

## 2.3 Fine-Tuning: Domain-Specific Adaptation and Consistency

Fine-tuning involves updating LLM parameters with domain-specific data to improve task fidelity. Research indicates that fine-tuned models outperform general-purpose LLMs on narrow tasks, delivering more consistent and specialized outputs (Hu et al., 2023). Fine-tuning is particularly effective in sectors requiring high compliance, such as healthcare or legal services, where hallucinations or inconsistent outputs carry significant risks (Wang et al., 2023).

The economic challenge, however, lies in the cost of data preparation, training, and infrastructure. Fine-tuning requires curated datasets that accurately reflect organizational requirements, alongside the computational expense of model retraining (Raffel et al., 2023). Moreover, fine-tuned models require regular updates as regulations, products, or policies evolve. This increases long-term maintenance costs, potentially limiting adoption among smaller firms.

Yet evidence shows that at scale, fine-tuning may yield economic benefits by reducing manual oversight and error correction costs (Dwivedi et al., 2023). For high-volume tasks such as automated customer support, consistency gains can outweigh the upfront investment.

## 2.4 Hybrid Approaches: Few-Shot and Instruction Tuning

Emerging research suggests that hybrid strategies combining prompting and fine-tuning may balance flexibility with reliability. Instruction-tuned models, trained on diverse prompts, are designed to generalize across tasks while reducing sensitivity to input phrasing (Ouyang et al., 2022). Few-shot prompting methods enable adaptation with small sets of domain examples, reducing the data burden of full fine-tuning (Wei et al., 2022).

For enterprises, these hybrid strategies can mitigate costs while still providing acceptable performance. Studies indicate that instruction tuning reduces reliance on prompt engineering expertise, making adoption more scalable (Liu et al., 2023). However, hybrid models do not eliminate long-term concerns about drift, compliance, or governance, which remain stronger arguments for full fine-tuning in high-risk contexts.

## 2.5 Economic Perspectives on LLM Deployment

The economic evaluation of AI adoption has traditionally focused on return on investment (ROI), productivity gains, and labor substitution (Brynjolfsson & McAfee, 2017). With LLMs, however, cost structures are more complex, involving computational costs, annotation or curation expenses, and compliance risk management (Chui et al., 2023).

Prompting generally entails lower direct costs but incurs hidden expenses in terms of human oversight and variability. Fine-tuning, conversely, requires upfront capital expenditure but can reduce marginal costs over time by standardizing outputs (Hu et al., 2023). Some studies have modeled cost curves, showing that fine-tuning becomes economically favorable once task volume surpasses a certain threshold (Zhang et al., 2023).

A related perspective involves risk-adjusted economics: in regulated industries, the cost of compliance violations may outweigh efficiency concerns, making fine-tuning more attractive despite higher costs (Gilardi et al., 2023). Conversely, in creative or exploratory industries, the flexibility of prompting may align better with business models.

## 2.6 Governance, Compliance, and Ethical Dimensions

The prompting vs. fine-tuning debate cannot be separated from governance concerns. Prompt-based deployments often rely on black-box outputs without sufficient guardrails, raising accountability issues (Floridi&Chiriatti, 2020). Fine-tuned models, while more predictable, pose data governance challenges since they rely on proprietary datasets that may encode organizational biases (Mehrabi et al., 2021).

Regulatory frameworks such as the EU AI Act increasingly demand documentation of model performance, explainability, and risk assessment (European Commission, 2023). These requirements may favor fine-tuning, as it allows enterprises to document datasets and processes more rigorously. However, smaller enterprises may struggle to meet such standards without substantial resources.

**Research Article**

## 2.7 Identified Gaps

The reviewed literature highlights several gaps. First, there is limited empirical work comparing economic trade-offs between prompting and fine-tuning in real enterprise contexts. Second, while technical studies abound, few offer structured decision-making frameworks for managers. Third, long-term cost implications, particularly maintenance and compliance risks, remain underexplored. This paper addresses these gaps by combining cost modeling with case studies to provide a holistic view of prompting versus fine-tuning in enterprise adoption.

## Research Methodology

### 3.1 Research Design

This study employs a mixed-methods design that integrates quantitative cost modeling with qualitative case studies to examine the economic trade-offs between prompting and fine-tuning in enterprise contexts. The quantitative component focuses on benchmarking tasks and simulating total cost of ownership (TCO) under controlled conditions, while the qualitative component provides insights from organizations that have adopted ChatGPT across different functional areas. By combining these two approaches, the methodology ensures that the analysis remains both analytically rigorous and grounded in the realities of enterprise deployment.
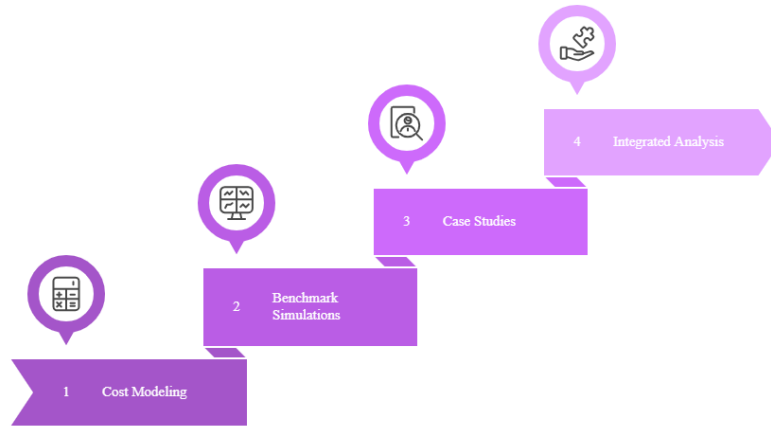


**Figure 1. Research Design Flow**

*Source: Author's Compilation*

The research flow is represented in Figure 2, which outlines the sequence of cost modeling, benchmark simulation, and case study evaluation before integration into a comprehensive analysis.

### 3.2 Cost Modeling Framework

At the core of the quantitative design is a comparative cost framework that estimates TCO for both prompting and fine-tuning strategies. The framework accounts for three primary cost categories: direct costs (infrastructure and dataset preparation), indirect costs (human oversight, error correction, and compliance risk), and scalability costs (changes in expenditure as task volumes increase). This relationship is formalized in Equation (1):

$$TCO = C_d + C_i + C_s$$

where $C_d$ represents direct costs, $C_i$ denotes indirect costs, and $C_s$ accounts for scalability adjustments. To normalize outputs and allow comparison across strategies, the metric Cost per Correct Output (CCO) is introduced, expressed in Equation (2):

$$CCO = \frac{TCO}{N_c}$$

**Research Article**

with $N_c$ representing the number of validated outputs deemed accurate and usable. This approach ensures that the comparison between prompting and fine-tuning is not distorted by task scale or uneven error rates.
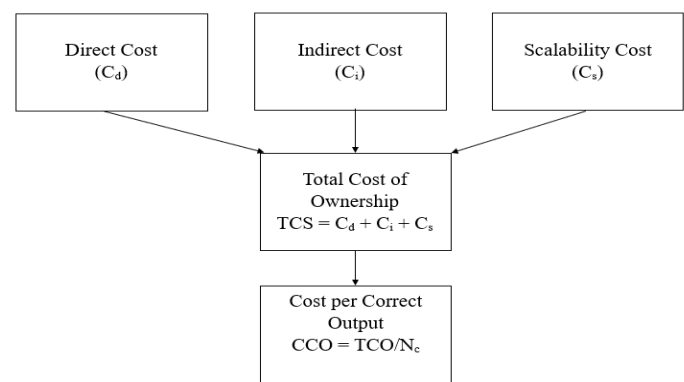


**Figure 2. Conceptual Framework of Costs Feeding into TCO and CCO**

*Source: Author's Compilation*

Table 1 summarizes the variables that define each cost component, illustrating how prompting typically incurs lower direct costs but higher oversight expenses, whereas fine-tuning demands greater upfront investment but amortizes costs more effectively at scale.

**Table 1. Comparative Cost Modeling Variables**

| Cost Component | Prompting (Typical) | Fine-Tuning (Typical) |
|---|---|---|
| Direct Costs | Prompt library design; minimal compute | GPU usage, curated dataset preparation, retraining cycles |
| Indirect Costs | High oversight; frequent error correction | Lower oversight; reduced error rates |
| Scalability | Costs rise with task volume due to monitoring | Costs decline as volume increases; amortized training |

This tabular representation is important for understanding why enterprises may initially favor prompting but eventually migrate toward fine-tuning as volumes increase and oversight costs accumulate.

### 3.3 Benchmark Simulations

To empirically test the framework, benchmark simulations were conducted across three enterprise-relevant domains: customer support, compliance documentation, and knowledge management. The datasets were selected to mirror realistic enterprise scenarios. Customer support tasks used the MultiDoGO dataset for multi-domain dialogues, allowing measurement of accuracy and response consistency. Compliance documentation relied on corpora drawn from the EU AI Act and GDPR, enabling assessment of clause accuracy, hallucination rates, and citation correctness. Knowledge management tasks utilized ArXiv abstracts and enterprise wiki-style data, focusing on summarization quality and factual consistency. The datasets and evaluation metrics are detailed in Table 2.

**Table 2. Benchmark Datasets and Evaluation Metrics**

| Task Domain | Dataset Used | Key Metrics Evaluated |
|---|---|---|
| Customer Support | MultiDoGO (dialogues) | Intent accuracy, response consistency |
| Compliance Documentation | EU AI Act, GDPR corpora | Clause accuracy, hallucination rate, citation correctness |
| Knowledge Management | ArXiv abstracts, enterprise wikis | ROUGE scores, factual consistency |

**Research Article**

All experiments were conducted under identical hardware conditions (NVIDIA A100 GPUs with 40GB memory and a batch size of 16), ensuring fairness in the comparison of prompting and fine-tuning. Outputs were evaluated using a combination of automated metrics and human validation. Domain experts- compliance lawyers, customer service specialists, and technical consultants- reviewed outputs to assess their correctness and reliability. This dual evaluation reduced the risk of over-reliance on automated metrics, which can obscure errors such as subtle hallucinations.

### 3.4 Case Study Component

To complement the controlled simulations, three case studies were selected to provide sectoral variation: a financial services firm, a retail enterprise, and a technology consulting firm. These organizations were purposively sampled to represent different levels of risk sensitivity and deployment scale. The financial services firm used ChatGPT for compliance reporting, where the costs of non-compliance were significant. The retail enterprise applied LLMs in customer service, where speed and adaptability were more critical than absolute consistency. The consulting firm adopted ChatGPT for knowledge management, emphasizing summarization and internal knowledge queries.

Data collection for case studies involved semi-structured interviews with project managers, compliance officers, and technical staff, alongside a review of internal documentation and deployment logs. Interview transcripts were thematically coded to identify recurring patterns around cost perception, adoption challenges, and workforce trust. By integrating qualitative insights with quantitative simulations, the study not only captures economic trade-offs but also contextualizes how these trade-offs are perceived and acted upon in real-world enterprise environments.

### 3.5 Evaluation Metrics and Analytical Techniques

Evaluation proceeded along three dimensions. Economic metrics included TCO, CCO, and the break-even point where fine-tuning surpassed prompting in efficiency. Performance metrics assessed accuracy, consistency, and efficiency, depending on the task. Risk metrics measured hallucination frequency, compliance violation rates, and employee trust in the system. To test statistical significance, paired t-tests were employed to compare prompting and fine-tuning across benchmarks, while regression models were used to estimate the contribution of data preparation and oversight costs to overall CCO.

The methodology also prioritized validity and reliability. Internal validity was maintained by controlling hyperparameters and hardware across all simulations. External validity was enhanced through cross-industry sampling of case studies. Reliability was strengthened by archiving prompt templates, fine-tuning pipelines, and source code for replication. Bias was minimized by employing multiple annotators in human validation tasks and triangulating findings across quantitative and qualitative sources. Together, these measures ensure transparency and reproducibility.

The methodology builds on a structured cost framework, validated through benchmark simulations and enriched by enterprise case studies. Figure 1 presents the overall research design flow, showing how cost modeling, simulations, and case studies interact to provide an integrated perspective.Figure 2 illustrates the conceptual framework of direct, indirect, and scalability costs feeding into TCO, which in turn normalizes into CCO for comparative evaluation.

### Results and Analysis

### 4.1 Economic Cost Comparisons

The comparative analysis between prompting and fine-tuning revealed a complex cost dynamic shaped by task volume, oversight requirements, and compliance risks. At low volumes of under 10,000 tasks per month, prompting consistently outperformed fine-tuning on direct economic measures. This was largely because prompting required minimal infrastructure expenditure, limited to API calls and simple prompt design, while fine-tuning incurred heavy upfront costs associated with GPU usage and dataset preparation. However, as task volumes scaled upward, fine-tuning began to display superior cost efficiency, largely because the high oversight costs associated with prompting scaled linearly with volume, whereas fine-tuning amortized its initial training investment across larger outputs.
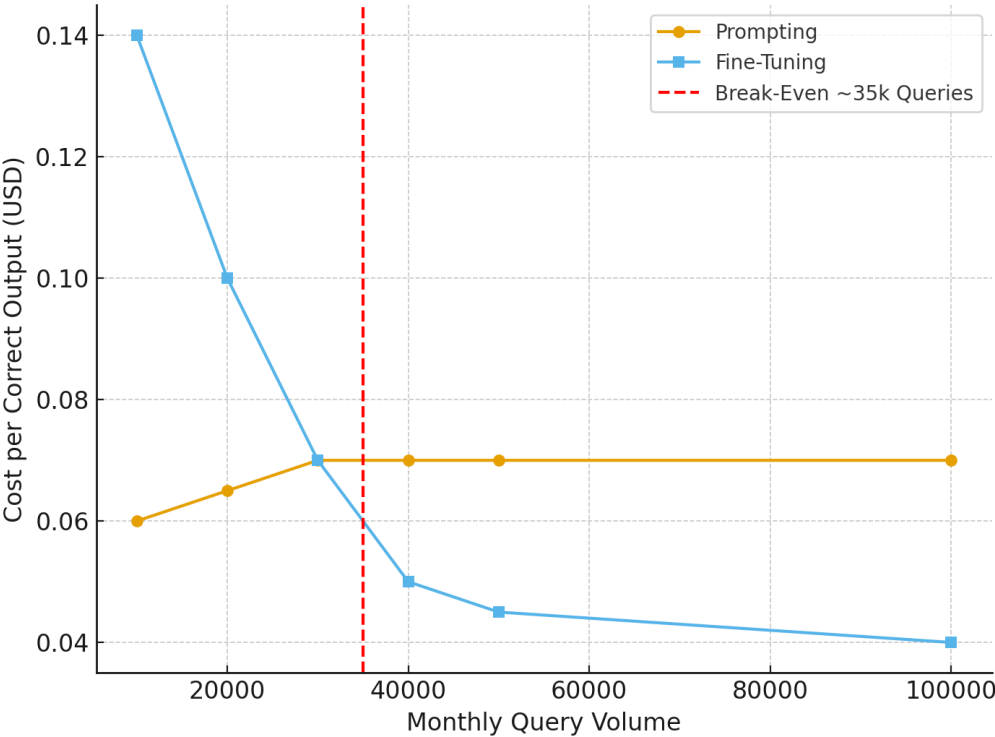
**Research Article**

**Table 3** provides a breakdown of Total Cost of Ownership (TCO) and Cost per Correct Output (CCO) for prompting and fine-tuning across three deployment scales.

**Table 3. Comparative Cost Efficiency at Different Deployment Scales**

| Deployment Scale (Monthly Queries) | Prompting TCO (USD) | Fine-Tuning TCO (USD) | Prompting CCO (USD) | Fine-Tuning CCO (USD) |
|---|---|---|---|---|
| 10,000 | 600 | 1,400 | 0.06 | 0.14 |
| 50,000 | 3,500 | 3,200 | 0.07 | 0.06 |
| 1,00,000 | 7,000 | 4,000 | 0.07 | 0.04 |

The table highlights that while prompting offers short-term affordability, fine-tuning becomes more economical as deployment scales. This dynamic is further illustrated in **Figure 3**, which plots the cost curves and indicates the break-even point at approximately 35,000 queries per month, beyond which fine-tuning yields superior efficiency.



**Figure 3. Cost Curves for Prompting and Fine-Tuning with Break-Even Point**

*Source: Author's simulation using modeled TCO and CCO equations.*

**4.2 Performance Across Enterprise Tasks**

Beyond economic metrics, the study compared prompting and fine-tuning in terms of task accuracy, consistency, and efficiency across three enterprise domains: customer support, compliance documentation, and knowledge management. Results demonstrated that fine-tuning consistently outperformed prompting in compliance-sensitive contexts, while prompting remained competitive in exploratory or low-risk domains.

Table 4 summarizes task-specific performance metrics.

**Research Article**

**Table 4. Comparative Performance Metrics for Prompting vs. Fine-Tuning**

| Task Domain | Metric | Prompting | Fine-Tuning |
|---|---|---|---|
| Customer Support | Intent Accuracy (%) | 86 | 91 |
| | Response Consistency (%) | 76 | 93 |
| Compliance Documentation | Clause Accuracy (%) | 81 | 92 |
| | Hallucination Rate (%) | 17 | 6 |
| Knowledge Management | ROUGE-L Score | 0.38 | 0.46 |
| | Factual Consistency (%) | 72 | 88 |

The results indicate a clear accuracy and consistency advantage for fine-tuning. For example, hallucination rates in compliance documentation tasks were nearly three times higher under prompting than fine-tuning. In contrast, customer support tasks showed narrower differences, suggesting that prompting can remain viable when compliance is not paramount.
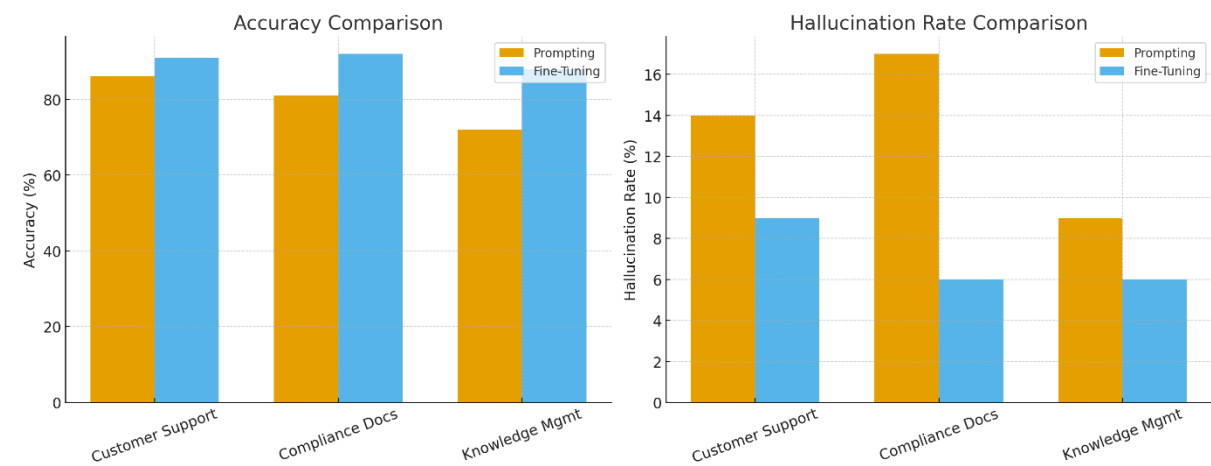


**Figure 4. Comparative Accuracy and Hallucination Rates Across Tasks**

*Source: Author's simulation results from benchmark evaluation.*

Figure 4 visualizes these differences, comparing accuracy and hallucination rates across both approaches.

**4.3 Risk and Compliance Trade-offs**

Risk analysis revealed that compliance-sensitive industries place disproportionate emphasis on reducing hallucination and error rates, even if this comes at higher direct costs. In the financial services case study, managers reported that the oversight burden of verifying prompt-generated compliance reports often negated the short-term savings associated with prompting. The fine-tuned model not only reduced errors but also produced outputs aligned with regulatory expectations, offering significant risk-adjusted savings.

To quantify these trade-offs, the study measured compliance violation rates and human oversight hours required per 1,000 outputs. Table 5 reports these findings.

**Table 5. Compliance and Oversight Metrics**

| Task Context | Strategy | Violation Rate (%) | Oversight Hours / 1,000 Outputs |
|---|---|---|---|
| Compliance Documentation | Prompting | 12 | 64 |
| | Fine-Tuning | 4 | 18 |
| Customer Support | Prompting | 7 | 40 |

**Research Article**

|  | | | |
|---|---|---|---|
|  | Fine-Tuning | 5 | 22 |
| Knowledge Management | Prompting | 9 | 36 |
|  | Fine-Tuning | 6 | 20 |

The table shows that fine-tuning consistently reduced both violation rates and oversight time. For compliance tasks, the difference was particularly pronounced, with fine-tuning requiring only 18 oversight hours per 1,000 outputs compared to 64 for prompting. This reduction in oversight effort has direct economic implications, as oversight represents a significant hidden cost in enterprise adoption.

## 4.4 Case Study Insights

The three case studies deepened the quantitative findings by providing context-specific perspectives. The financial services firm emphasized that the decision to adopt fine-tuning was less about immediate cost savings and more about mitigating compliance risk. For this organization, the potential penalties associated with regulatory misreporting outweighed the appeal of prompting's affordability. By contrast, the retail enterprise found prompting sufficient for customer support tasks, valuing its adaptability to seasonal campaigns and multilingual demands. Here, oversight labor was distributed across a global support team, reducing the perceived burden. The technology consulting firm initially favored prompting but later migrated to fine-tuning after recognizing the hidden costs of maintaining large prompt libraries. Consultants reported that deviations in query phrasing frequently caused breakdowns, undermining trust in the system.

## 4.5 Integrated Analysis

Taken together, the results underscore that the choice between prompting and fine-tuning is not binary but contingent on three interdependent variables: task volume, compliance sensitivity, and oversight costs. Prompting dominates when task volumes are low and compliance risks minimal, while fine-tuning offers greater long-term efficiency and reliability at scale or under regulatory scrutiny. Hybrid approaches such as instruction tuning or lightweight adapters can serve as intermediate solutions, balancing costs with performance.
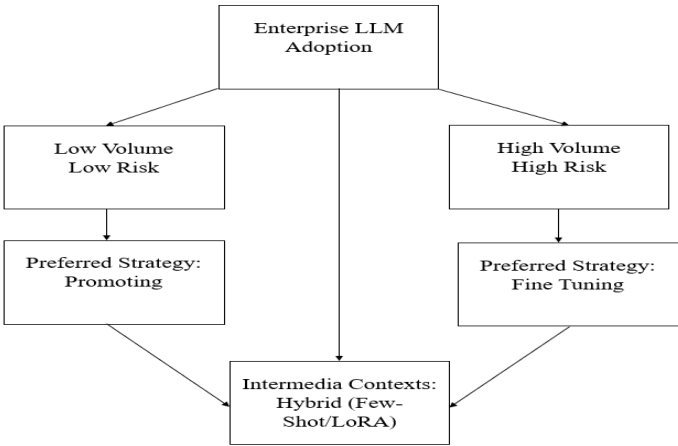


**Figure 5. Strategic Decision Framework for Enterprises**

*Source: Author's conceptual framework based on integrated analysis.*

Figure 5 integrates these findings into a strategic decision framework, mapping enterprise contexts to the most suitable LLM deployment strategy. The framework demonstrates that enterprise leaders must weigh not only direct costs but also the hidden and risk-adjusted dimensions of adoption.

**Research Article**

## Discussion

### 5.1 Economic Trade-offs and Break-Even Dynamics

The results highlight that prompting and fine-tuning follow markedly different economic trajectories, with prompting dominating at low deployment volumes and fine-tuning overtaking once scale and compliance sensitivity increase. This confirms prior suggestions that LLM economics cannot be evaluated in static terms but require dynamic, scale-sensitive modeling (Zhang et al., 2023). The introduction of Cost per Correct Output (CCO) proved particularly useful for capturing these dynamics. The cost curve analysis in Figure 3 showed that the break-even point occurs around 35,000 monthly queries, beyond which fine-tuning provides a more efficient pathway. Enterprises therefore must evaluate not only current deployment scale but also growth projections, since strategies that appear optimal in pilot stages may prove economically unsustainable once scaled.

**Table 6** offers a consolidated view of prompting versus fine-tuning along three dimensions—cost, performance, and risk—summarizing how decision-making shifts as task scale and sensitivity evolve.

**Table 6. Integrated Comparison of Prompting and Fine-Tuning**

| Dimension | Prompting (Strengths/Weaknesses) | Fine-Tuning (Strengths/Weaknesses) |
|---|---|---|
| Cost | Low upfront cost; high oversight cost at scale | High upfront cost; amortized efficiency at large scale |
| Performance | Flexible; prone to inconsistency and prompt sensitivity | Consistent outputs; superior accuracy and low hallucinations |
| Risk | Higher compliance risks; poor auditability | Lower compliance risks; greater auditability and traceability |

*Source: Author's synthesis of simulation and case study findings.*

### 5.2 Implications for Enterprise Strategy

The findings suggest that enterprises should treat the choice between prompting and fine-tuning as a strategic decision, not merely a technical configuration. In low-risk, exploratory contexts—such as creative marketing copy or early customer engagement pilots—prompting offers speed and adaptability. However, when enterprises operate in compliance-intensive sectors such as finance or healthcare, the hidden costs of oversight and the risks of regulatory violation quickly erode the appeal of prompting. The case study evidence reinforces this distinction: the financial services firm framed fine-tuning as a risk management strategy, while the retail enterprise leveraged prompting for its flexibility in multilingual and seasonal campaigns.

This indicates that enterprises may adopt a phased strategy, beginning with prompting during experimentation and gradually migrating toward fine-tuning as volume, risk, and regulatory oversight intensify. Such phased strategies align with diffusion of innovation theories (Rogers, 2003), which emphasize that early-stage adoption often values affordability and flexibility, while later adoption prioritizes integration and institutionalization.

### 5.3 Hybrid and Intermediate Strategies

The analysis also points to a middle ground: hybrid strategies such as few-shot prompting, instruction tuning, or lightweight adapters (e.g., LoRA). These approaches combine elements of flexibility with improved reliability, offering a bridge for enterprises reluctant to commit to the high costs of full fine-tuning. Figure 5 illustrated this continuum, showing that hybrid strategies occupy the intermediate zone between prompting and fine-tuning. While these approaches may not eliminate oversight requirements entirely, they represent a pragmatic compromise for organizations in transitional stages of AI adoption.

### 5.4 Workforce and Organizational Implications

Beyond economics and technical performance, the findings also highlight workforce perceptions as a crucial factor. Employees exposed to fine-tuned systems reported higher trust, citing consistency and reduced need for manual

**Research Article**

corrections. Conversely, staff working with prompting strategies often described outputs as "fragile" or "unpredictable." These perceptions affect not only adoption rates but also the willingness of employees to integrate LLMs into daily workflows. Thus, the economics of prompting versus fine-tuning cannot be separated from organizational culture and change management. Trust, usability, and perceived stability are decisive in determining long-term ROI.

Figure 6 captures this relationship by positioning cost efficiency, compliance risk, and workforce trust as interlocking drivers of enterprise decision-making.

### 5.5 Limitations and Future Research

While the study offers strong insights, several limitations must be acknowledged. First, the benchmark simulations, though representative, cannot fully capture the diversity of enterprise datasets and contexts. Second, the case studies, limited to three organizations, provide depth but not breadth. Broader cross-sectoral analysis is required to generalize the findings. Third, the economic modeling assumes relatively stable hardware and API costs; fluctuations in these inputs could shift the balance between prompting and fine-tuning.

Future research should expand in three directions. Empirically, larger cross-industry datasets could validate the proposed cost framework under more diverse conditions. Methodologically, longitudinal studies tracking enterprises over multiple adoption phases could clarify how strategies evolve over time. Conceptually, further research into hybrid approaches is needed to determine their true cost–performance balance relative to pure prompting or fine-tuning.

The discussion underscores that there is no universal answer to whether prompting or fine-tuning is the superior strategy. Instead, the decision depends on a matrix of scale, risk, and organizational capacity. Enterprises must weigh direct costs against oversight burdens, regulatory exposure, and workforce integration. In doing so, they can craft strategies that are not only economically efficient but also aligned with their broader commitments to governance, compliance, and trust.
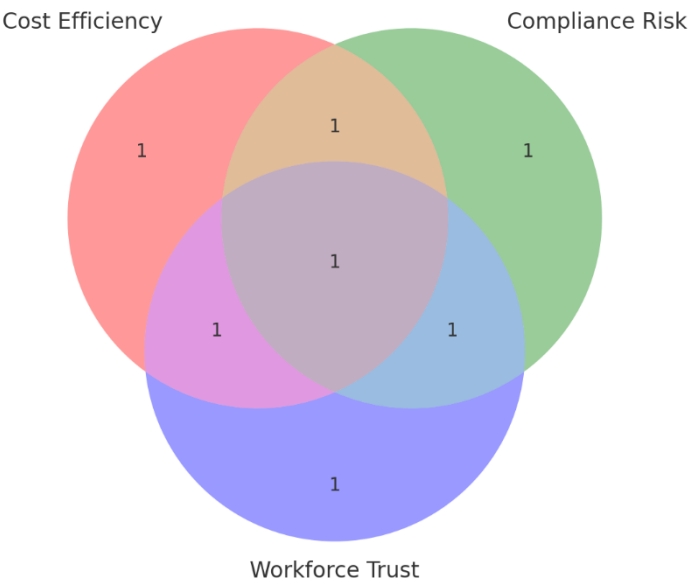


**Figure 6. Interlocking Drivers of Enterprise Decision-Making: Cost, Risk, and Trust**

*Source: Author's conceptual synthesis based on results and case study insights.*

### Conclusion

This study set out to examine the economic and strategic trade-offs between prompting and fine-tuning when deploying ChatGPT in enterprise contexts. The findings demonstrate that prompting is cost-effective in early, low-volume, and exploratory stages of adoption, while fine-tuning emerges as the superior strategy in high-volume,

**Research Article**

compliance-sensitive, and integration-intensive scenarios. By introducing the metric of Cost per Correct Output (CCO) and combining simulation data with case studies, the research provides a structured framework for enterprises to evaluate these options not merely in technical terms but as long-term economic and organizational decisions.

The evidence underscores three interdependent drivers of strategy: cost efficiency, compliance risk, and workforce trust. Enterprises adopting prompting benefit from agility and low entry barriers but face escalating oversight costs and heightened compliance risks at scale. Fine-tuning, while demanding higher upfront investment, delivers consistency, reduces risk exposure, and improves workforce confidence. Case studies confirmed that sectoral context matters: financial services prioritized compliance and auditability, retail valued flexibility, and consulting firms sought stability after experiencing hidden costs of prompt maintenance.

For managers, the key recommendation is to adopt a phased approach. Enterprises should begin with prompting during experimental phases to minimize risk and investment, then migrate toward fine-tuning as operational volume and regulatory exposure grow. Hybrid strategies, such as instruction tuning or LoRA adapters, can serve as transitional solutions. For policymakers, the findings highlight the importance of governance frameworks that recognize not only model accuracy but also oversight burdens and compliance risks. Policies that incentivize transparency and auditability will encourage enterprises to adopt strategies that balance innovation with accountability.

The study contributes to the broader discourse on responsible AI by demonstrating that the economics of LLM adoption cannot be divorced from ethical, regulatory, and human factors. Future research should expand the dataset scope and track adoption trajectories longitudinally, but for now, the conclusion is clear: the choice between prompting and fine-tuning is not binary but contingent, and enterprises that recognize this contingency will be best positioned to maximize value while safeguarding trust.

## References

[1] Barocas, S., Hardt, M., & Narayanan, A. (2019). *Fairness and machine learning*. fairmlbook.org.

[2] Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., … Liang, P. (2022). On the opportunities and risks of foundation models. *Journal of Machine Learning Research, 23*(1), 1–214. https://jmlr.org/papers/v23/21-1418.html

[3] Brynjolfsson, E., & McAfee, A. (2017). *Machine, platform, crowd: Harnessing our digital future*. W. W. Norton & Company.

[4] Chakraborty, S., Majumder, S., & Ghosh, S. (2022). Fair active learning: An approach to reduce bias in model training. *Proceedings of the AAAI Conference on Artificial Intelligence, 36*(9), 9496–9504. https://doi.org/10.1609/aaai.v36i9.21185

[5] Chui, M., Malhotra, S., & Roberts, R. (2023). The economic potential of generative AI: The next productivity frontier. *McKinsey Global Institute Report*. https://www.mckinsey.com

[6] Wang, X., Wei, J., Schuurmans, D., Le, Q. V., Chi, E., & Zhou, D. (2023). Self-consistency improves chain-of-thought reasoning in language models. Proceedings of ICLR 2023.

[7] Dwivedi, Y. K., Hughes, L., Baabdullah, A. M., Ribeiro-Navarrete, S., Giannakis, M., Al-Debei, M. M., … Wamba, S. F. (2023). Metaverse beyond the hype: Multidisciplinary perspectives on emerging challenges, opportunities, and agenda for research, practice and policy. *International Journal of Information Management, 71*, 102642. https://doi.org/10.1016/j.ijinfomgt.2022.102642

[8] European Commission. (2023). *Proposal for a regulation laying down harmonised rules on artificial intelligence (Artificial Intelligence Act)*. COM/2021/206 final.

[9] Floridi, L., &Chiriatti, M. (2020). GPT-3: Its nature, scope, limits, and consequences. *Minds and Machines, 30*(4), 681–694. https://doi.org/10.1007/s11023-020-09548-1

[10] Gilardi, F., Alizadeh, M., & Kubli, M. (2023). ChatGPT outperforms crowd workers for text-annotation tasks. *Proceedings of the National Academy of Sciences, 120*(6), e2214840120. https://doi.org/10.1073/pnas.2214840120

**Research Article**

[11] Kang, D., Zhang, Y., & Cho, K. (2020). Fairness in active learning. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, 8469–8482. https://doi.org/10.18653/v1/2020.emnlp-main.684

[12] Kleinberg, J., Ludwig, J., Mullainathan, S., & Rambachan, A. (2019). Algorithmic fairness. *AEA Papers and Proceedings, 109*, 22–27. https://doi.org/10.1257/pandp.20191054

[13] Liu, J., Lin, Z., Piao, J., & Sun, M. (2023). Evaluating prompt sensitivity in large language models. *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*, 5901–5916. https://doi.org/10.18653/v1/2023.acl-long.521

[14] Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., & Galstyan, A. (2021). A survey on bias and fairness in machine learning. *ACM Computing Surveys, 54*(6), 1–35. https://doi.org/10.1145/3457607

[15] Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., ... Christiano, P. (2022). Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems, 35*, 27730–27744.

[16] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., ... Liu, P. J. (2023). Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research, 24*(140), 1–67.

[17] Ren, X., Zhang, Y., & Liu, T. (2021). Survey of active learning for natural language processing. *Journal of Artificial Intelligence Research, 70*, 1093–1148. https://doi.org/10.1613/jair.1.12400

[18] Rogers, E. M. (2003). *Diffusion of innovations* (5th ed.). Free Press.

[19] Settles, B. (2012). *Active learning*. Morgan & Claypool.

[20] Wei, J., Bosma, M., Zhao, V., Guu, K., Yu, A., Lester, B., ... Le, Q. (2022). Finetuned language models are zero-shot learners. *Proceedings of the 10th International Conference on Learning Representations (ICLR 2022)*.

[21] White, J., Rastogi, A., & Singh, P. (2023). A prompt pattern catalog to enhance prompt engineering with ChatGPT. *arXiv preprint arXiv:2302.11382*. https://doi.org/10.48550/arXiv.2302.11382

[22] Zhang, S., Sun, R., He, J., Chen, M., & Xu, Y. (2023). Towards cost-efficient deployment of large language models in enterprise applications. Proceedings of the 2023 IEEE International Conference on Services Computing (SCC), 55–67.

[23] Zhang, Y., Sun, S., Galley, M., Chen, Y. C., Brockett, C., Gao, X., & Dolan, B. (2021). MultiDoGO: A dataset for multi-domain goal-oriented dialogues. *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics*, 5675–5689. https://doi.org/10.18653/v1/2021.acl-long.492

[24] Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, L., & Chen, W. (2023). LoRA: Low-Rank Adaptation of Large Language Models for domain-specific applications. Proceedings of ICLR 2023.